

# 基于 ECOC 的多类代价敏感分类方法

吴崇明<sup>1</sup> 王晓丹<sup>2</sup> 薛爱军<sup>2</sup> 来杰<sup>2</sup>

1 西京学院商学院 西安 710123

2 空军工程大学防空反导学院 西安 710051

(w\_9887@163.com)

**摘要** 研究了基于纠错输出编码实现多类代价敏感分类的方法,提出了一种新的将多类代价敏感分类问题分解为多个二类代价敏感分类问题的框架。为获得其中每个二类代价敏感基分类器的二类代价矩阵,提出了利用已知多类代价矩阵计算误分类代价的期望值的方法,给出了计算二类代价矩阵的通用计算公式。为验证所提方法的有效性,在人工和 UCI 数据集上将其与现有方法进行了比较,实验结果表明所提方法具有相似甚至更好的性能。

**关键词:** 多类代价敏感分类;纠错输出编码;多类代价矩阵;二类代价矩阵

中图分类号 TP391

## Multiclass Cost-sensitive Classification Based on Error Correcting Output Codes

WU Chong-ming<sup>1</sup>, WANG Xiao-dan<sup>2</sup>, XUE Ai-jun<sup>2</sup> and LAI Jie<sup>2</sup>

1 Business School, Xijing University, Xi'an 710123, China

2 College of Air and Missile Defense, Air force Engineering University, Xi'an 710051, China

**Abstract** Approach of multiclass cost-sensitive classification based on error correcting output codes is studied in this paper, and a new framework to decompose the complex multiclass cost-sensitive classification problem into a series of binary cost-sensitive classification problems is proposed. In order to obtain the binary cost matrix of each binary cost-sensitive base classifier, a method of computing the expected misclassification costs from the given multiclass cost matrix is proposed, and the general formula for computing the binary costs are given. Experimental results on artificial datasets and UCI datasets show that the proposed method has similar or even better performance in comparison with the existing methods.

**Keywords** Multiclass cost-sensitive classification, Error correcting output codes, Multiclass cost matrix, Binary cost matrix

## 1 引言

对许多实际应用问题,如疾病诊断<sup>[1]</sup>和欺诈检测<sup>[2]</sup>等,不同类型的分类错误造成的误分代价明显不同,这些都属于代价敏感分类问题。为解决代价敏感分类问题,研究者提出了多种方法<sup>[3-9]</sup>。根据实现策略的不同,代价敏感学习主要分为两大类:间接代价敏感学习(或代价敏感元学习)<sup>[1,4]</sup>和直接代价敏感学习<sup>[5-7]</sup>。间接代价敏感学习并不改变传统的学习算法,而是将代价敏感学习方法作为一个单独的模块,对训练数据进行预处理或者对学习算法的输出结果进行后处理,从而实现代价敏感学习。直接代价敏感学习是直接对分类器引入到分类器的构建中,通过更改分类器的结构或者修改分类器的目标函数,将传统的以最小错误率为目标的分类器转换为以最小风险为目标,从而实现代价敏感学习。

对于如何实现多类代价敏感分类,将纠错输出编码(Error-Correcting Output Codes, ECOC)<sup>[10-12]</sup>多类分解框架与代价敏感二分类器结合引起了研究者的关注。文献<sup>[14]</sup>将多类代价敏感分类问题分解为一系列二分类代价敏感分类问题,通过集成方法由各二分类代价敏感基分类器的输出生成多类

分类结果。Lin<sup>[15]</sup>利用代价转换技术将“一对一”和“一对多”编码方法扩展到了多类代价敏感分类问题。与 Lin 的工作相比,Langford<sup>[16]</sup>提出了一个更加通用的框架,该框架基于修改后的 ECOC,将多类代价敏感分类问题分解为若干个二类代价敏感分类问题。但是,该框架的性能较差,其中包含的取阈值的步骤,使得代价信息的表达变得不准确<sup>[17]</sup>。如何设计更有效、更通用的多类代价敏感分类方法仍是具有挑战性的问题。

基于 ECOC 多类分解框架实现多类代价敏感分类方法,通过 ECOC 多类分解框架将多类代价敏感分类问题分解为一系列二类代价敏感分类问题,存在的一个关键问题:如何从给定的多类代价矩阵出发得到每个类的错分代价,并设计新的、更加有效的分解框架;如何在将多类代价敏感分类问题分解为一系列二类分类代价敏感问题的同时,由给定的多类代价矩阵得出对应的一系列二类代价矩阵。本文基于原始 ECOC 框架,提出了一种新的分解框架;为获得其中每个二类代价敏感基分类器的二类代价矩阵,提出了利用已知多类代价矩阵计算误分类代价的期望值的方法;最后给出了计算二类代价矩阵的通用计算公式。

基金项目:国家自然科学基金(61876189,61273275,61703426)

This work was supported by the National Natural Science Foundation of China(61876189,61273275,61703426).

通信作者:王晓丹(afeu\_w@163.com)

## 2 二类代价敏感纠错输出编码

纠错输出编码基于二元{1, -1}或三元{1, 0, -1}的编码矩阵来实现类别的分解和基分类器输出结果的融合。{1}和{-1}分别代表一类子类, {0}代表对应的类别不参与相应基分类器的生成。图1给出了一个4类样本数据的编码矩阵示意图。

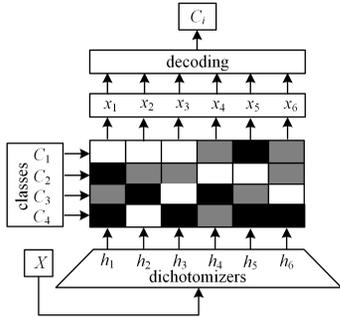


图1 4类模式 ECOC 分类

Fig.1 4 classes ECOC classification

编码矩阵中行代表类别, 列代表基分类器。原始类别通过每一列对应的二类划分成正、负子类, 分别用1和-1表示, 即白色和黑色, 灰色代表码字0。在获得二类划分之后, 用新的正负子集训练基分类器, 并对待识别样本进行测试, 输出结果, 得到输出向量(结果可以为软输出也可以为硬输出), 最后确定解码策略进行基分类器输出结果的融合。

常见的“一对一”和“一对多”多类分类方法属于 ECOC 分解框架的两种特殊情况, 即使是层次分解、决策树分解, 也能得到对应的 ECOC 编码矩阵。ECOC 方法简单有效, 避免了对复杂多类分类问题直接进行建模, 使多类分类能够利用二分类中经典的成熟理论和方法, 同时继承了纠错输出编码的容错能力<sup>[10]</sup>。

对于代价敏感多类分类问题, 多类代价矩阵 (Multiclass Cost Matrix, MCM) 表示了各类不同的误分类的代价, 例如, 4类代价矩阵如式(1)所示:

$$\begin{bmatrix} \omega_{11} & \omega_{12} & \omega_{13} & \omega_{14} \\ \omega_{21} & \omega_{22} & \omega_{23} & \omega_{24} \\ \omega_{31} & \omega_{32} & \omega_{33} & \omega_{34} \\ \omega_{41} & \omega_{42} & \omega_{43} & \omega_{44} \end{bmatrix} \quad (1)$$

其中,  $\omega_{ij}$  表示将一个属于  $c_i$  类的样本错分为属于  $c_j$  类的代价或风险。

对于代价敏感多类分类问题, 也可按照 ECOC 多类分解框架实现, 但需要同时考虑多类代价矩阵, 以及分解后每个二类代价敏感基分类器的二类代价矩阵, 设计分类器的目的是最小化误分类代价。基于上述考虑, 本文给出一个新的 ECOC 分解框架, 称之为二类代价敏感纠错输出编码 (Binary Cost Sensitive Error Correcting Output Codes, BCS\_ECOC)。图2给出了4类 BCS\_ECOC 的一个示例。

BCS\_ECOC 编码矩阵的定义与 ECOC 中编码矩阵的定义相同。同时, BCS\_ECOC 也包含两个重要的阶段: 训练阶段和测试阶段。在训练阶段, 训练数据集被重组为两个超类: 正类和负类。这一过程根据编码矩阵每一列元素的符号号完成。然后, 基于各二类分类基分类器的代价矩阵, 二类代价敏

感基分类器在重组后的训练数据集上完成训练。在测试阶段, 利用训练得到的代价敏感基分类器对未知类别的样本进行分类, 得到的分类结果将根据解码规则进行解码, 这与 ECOC 中的相同。

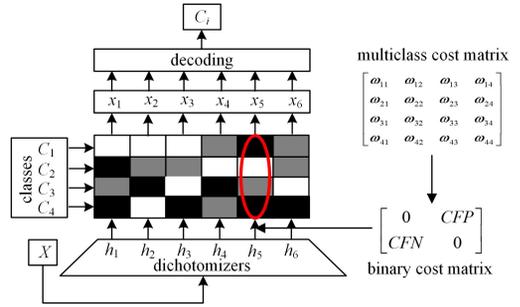


图2 BCS\_ECOC 的一个示例

Fig.2 Example of BCS\_ECOC

本文二类代价敏感纠错输出编码 BCS\_ECOC 与传统纠错输出编码 ECOC 的不同体现在: 1) 设计目标不同, ECOC 的设计目标是 minimized 错误分类率, BCS\_ECOC 的设计目标是 minimized 误分类代价; 2) 基分类器不同, 传统 ECOC 中的基分类器是非代价敏感的二类分类器, 如支持向量机, BCS\_ECOC 中的基分类器则是代价敏感的二类分类器, 如代价敏感支持向量机。基于 BCS\_ECOC, 许多成熟的二类代价敏感分类器, 都可被扩展用于解决多类代价敏感分类问题。

在上述 BCS\_ECOC 框架中, 需要解决的一个关键问题是: 如何基于已知的多类代价矩阵为每个二类代价敏感基分类器计算对应的二类代价矩阵。下节将给出解决。

## 3 基分类器二类代价矩阵的生成

如何根据已知的多类代价矩阵为每个二类代价敏感基分类器计算二类代价矩阵, 是本节要解决的问题。以下首先探讨当编码矩阵为“一对一”编码时二类代价矩阵的计算方法, 然后探讨当编码矩阵为“一对多”编码时二类代价矩阵的计算方法, 最后给出基于多类代价矩阵计算二类代价矩阵的通用计算公式。

### 3.1 “一对一”编码的二类代价矩阵计算

当编码矩阵为“一对一”时, 为每个二类代价敏感基分类器计算二类代价矩阵相对比较容易, 因为这种情况下(训练基分类器的)正类和负类都只包含一个类。图3给出了一个“一对一”编码的例子, 其中包含的类别个数为4。在图3中, 以基分类器  $f_1$  为例, 其正类 (Positive Class) 为  $c_1$ , 负类 (Negative Class) 为  $c_2$ ; 对于二类代价敏感基分类器  $f_1$ , 需要确定如式(2)所示的代价矩阵。

	$f_1$	$f_2$	$f_3$	$f_4$	$f_5$	$f_6$
Positive Class →	$c_1$	+1	+1	0	0	0
Negative Class →	$c_2$	-1	0	0	+1	0
	$c_3$	0	-1	0	-1	0
	$c_4$	0	0	-1	0	-1

图3 “一对一”编码某一行对应的正类和负类

Fig.3 Two classes in one of dichotomizers induced by one-to-one coding matrix

$$P \begin{matrix} & N \\ 0 & CFN \\ N \begin{bmatrix} CFP & 0 \end{bmatrix} \end{matrix} \quad (2)$$

其中,CFN 为错分为负类的代价(Cost for a False Negative, CFN),即正类误分类的代价,此处为属于正类的样本被错分为属于负类的代价;CFP 为错分为正类的代价(Cost for a False Positive, CFP),即负类误分类的代价,此处为属于负类的样本被错分为属于正类的代价。对于要分类  $c_1$  和  $c_2$  这两个类的基分类器  $f_1$ ,其代价矩阵中的这两个错分代价 CFN 和 CFP 都可以从给定的多类代价矩阵中直接获得。

CFN 就是将一个属于类  $c_1$  的样本错分为属于类  $c_2$  的代价。在如式(1)所示的多类代价矩阵中,CFN= $\omega_{12}$ 。另一方面,CFP 是将一个属于类  $c_2$  的样本错分为属于类  $c_1$  的代价,在多类代价矩阵中,CFP= $\omega_{21}$ 。所以,基分类器  $f_1$  的二类代价矩阵如式(3)所示。同理,可以得到其他基分类器的二类代价矩阵。

$$P \begin{matrix} & N \\ \omega_{11} & \omega_{12} \\ N \begin{bmatrix} \omega_{21} & \omega_{22} \end{bmatrix} \end{matrix} \quad (3)$$

### 3.2 “一对多”及其他编码方法的二类代价矩阵计算

对于“一对多”及其他编码矩阵来说,计算某一个基分类器的二类代价矩阵相对比较困难。其原因是,此时在二类代价敏感分类问题中,正类和负类通常会包含多个原始类。图 4 给出了一个“一对多”编码矩阵的例子,它包含的类别数为 4。

	$f_1$	$f_2$	$f_3$	$f_4$
Positive Class →	$c_1$ +1	-1	-1	-1
	$c_2$ -1	+1	-1	-1
Negative Class →	$c_3$ -1	-1	+1	-1
	$c_4$ -1	-1	-1	+1

图 4 “一对多”编码中基分类器对应的正类和负类

Fig. 4 Two classes in one of dichotomizers induced by one-to-many coding matrix

在图 4 中,仍然以基分类器  $f_1$  作为示例。可以看出,对于  $f_1$ ,正类仅由类  $c_1$  组成,而负类由 3 个类  $c_2, c_3$  和  $c_4$  组成。此时,错分为负类的代价 CFN 包含了 3 个方面:将属于类  $c_1$  的样本错分为属于类  $c_2$  的代价,将属于类  $c_1$  的样本错分为属于类  $c_3$  的代价,将属于类  $c_1$  的样本错分为属于类  $c_4$  的代价。另一方面,错分为正类的代价 CFP 同样包含 3 个方面:将属于类  $c_2$  的样本错分为属于类  $c_1$  的代价,将属于类  $c_3$  的样本错分为属于类  $c_1$  的代价,将属于类  $c_4$  的样本错分为属于类  $c_1$  的代价。计算基分类器  $f_1$  的二类代价矩阵,关键在于如何建立 CFN 和 CFP 与这 3 个方面的代价之间的联系。

我们认为,正类误分类的代价 CFN 和负类误分类的代价 CFP 可以通过计算误分类代价的期望值得到。假设正类包含的类的集合为  $S_p$ ,负类包含的类的集合为  $S_N$ 。对于如图 4 所示的基分类器  $f_1$ ,有  $S_p = \{c_1\}$  和  $S_N = \{c_2, c_3, c_4\}$ 。以 CFN 的计算为例进行说明,正类误分类的代价 CFN 的计算公式为:

$$CFN = \sum_{c_i \in S_p} P(c_i) \sum_{c_j \in S_N} P(c_i, c_j) \omega(c_i, c_j) \quad (4)$$

其中, $P(c_i)$ 是类  $c_i$  的先验概率,它满足条件  $\sum_{c_i \in S_p} P(c_i) = 1$ 。 $P(c_i, c_j)$ 是属于类  $c_i$  的样本错分为属于类  $c_j$  的概率, $\omega(c_i,$

$c_j)$ 是属于类  $c_i$  的样本错分为属于类  $c_j$  的代价,它的值可以从多类代价矩阵中直接找到。特别地,需要注意的是概率  $P(c_i, c_j)$ 必须满足条件  $\sum_{c_j \in S_N} P(c_i, c_j) = 1, c_i \in S_p$ 。

同理,负类误分类的代价 CFP 的计算公式为:

$$CFP = \sum_{c_i \in S_N} P(c_i) \sum_{c_j \in S_p} P(c_i, c_j) \omega(c_i, c_j) \quad (5)$$

在式(4)和式(5)中,某一类的先验概率  $P(c_i)$ 可以根据每一类包含样本的个数估计得到,错分代价  $\omega(c_i, c_j)$ 可以在多类代价矩阵中找到。计算 CFN 和 CFP 的困难在于如何得到错分概率  $P(c_i, c_j)$ ,可以用两类之间的错误分类率近似估计错分概率。

通常,如果两类之间的类别的可分性越好,则分类错误率就越低,因此,错分概率  $P(c_i, c_j)$ 与类别可分性成反比关系。所以,错分概率  $P(c_i, c_j)$ 的计算公式可表示为:

$$P(c_i, c_j) = \frac{1}{J(c_i, c_j)} \quad (6)$$

其中, $J(\cdot, \cdot)$ 为类别可分性的度量准则。常见的类别可分性度量有类内类间距离、Bhattacharyya 距离、Chernoff 界等。

需要注意的是,式(4)和式(5)可以用来为不同种类的编码矩阵计算二类代价矩阵,比如,密集随机编码、稀疏随机编码以及其他基于问题的编码方法等。因此,这两个公式可以作为根据多类代价矩阵计算二类代价矩阵的通用公式。

以下通过实验验证这两个公式的有效性,同时检验 BCS\_ECOC 的分类性能。

## 4 实验及分析

本节从实验数据、实验设计以及实验结果与分析 3 个方面对实验进行介绍。

### 4.1 实验数据

实验选取了两种不同的数据集,分别是人工数据集和 UCI 数据集。

#### (1)人工数据集

应用人工数据集的主要优势在于,人工数据集中每一类的条件概率密度函数是已知的,可以更加精确地计算类别之间的可分性。例如,基于每个类已知的条件概率可以计算两类之间的 Bhattacharyya 距离,它可以比类内类间距离更加准确地反映类间的可分性。这些人工数据集共包含 6 类,每一类含有相同数目的样本,其样本分布如图 5 所示,每一类的特征向量的维数为 2,分别用 Feature1 和 Feature2 来表示。每一类的特征向量都服从正态分布,其概率密度函数定义如式(7)所示,其中的参数设置见表 1。

$$p(x|class_i) = \frac{1}{2\pi|\Sigma_i|^{\frac{1}{2}}} \exp \left[ -\frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) \right] \quad (7)$$

$i = 1, 2, \dots, 6$

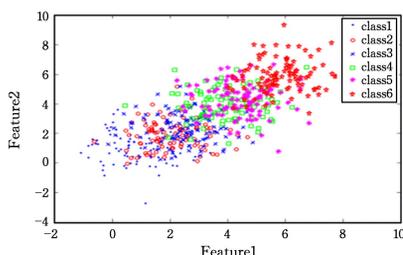


图 5 人工数据集中的数据分布

Fig. 5 Data distribution in artificial datasets

表1 人工数据集的参数设置

Table 1 Parameters settings for artificial datasets

Class	Prior probabilities	Mean vectors	Covariance matrices
$C_1$	$P(C_1) = \frac{1}{6}$	$\mu_1 = (1, 1)^T$	$\Sigma_1 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$
$C_2$	$P(C_2) = \frac{1}{6}$	$\mu_2 = (2, 2)^T$	$\Sigma_1 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$
$C_3$	$P(C_3) = \frac{1}{6}$	$\mu_3 = (2.5, 2.5)^T$	$\Sigma_3 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$
$C_4$	$P(C_4) = \frac{1}{6}$	$\mu_4 = (4, 4)^T$	$\Sigma_4 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$
$C_5$	$P(C_5) = \frac{1}{6}$	$\mu_5 = (4.5, 4.5)^T$	$\Sigma_2 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$
$C_6$	$P(C_6) = \frac{1}{6}$	$\mu_6 = (6, 6)^T$	$\Sigma_2 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$

在相同的参数下,为了进一步增加实验的鲁棒性,共生成了6个具有不同样本数目的人工数据集。这6个人工数据集中每类包含的样本的数目分别为10, 50, 100, 200, 300和400。为了表述方便,在下文中分别用1#, 2#, 3#, 4#, 5#和6#来代表这6个人工数据集。

## (2)UCI数据集

UCI数据集<sup>[18]</sup>常被用来衡量分类方法的性能,实验从UCI数据集中选取了14个多类的数据集。表2列出了这些数据集的样本的数目、特征量的个数以及包含类别的个数。

表2 实验中用到的UCI数据集

Table 2 UCI data sets used in experiment

Dataset	Samples	Attributes	Classes
balance	625	4	3
glass	214	10	6
iris	150	4	3
letter	1214	16	26
page-blocks	5473	10	5
sat	6435	36	6
segmentation	2310	19	7
soybean	306	35	18
thyroid	215	5	3
vehicle	846	18	4
vowel	990	13	11
wine	178	13	3
yeast	1484	8	10
zoo	101	16	7

## 4.2 实验设计

### (1)实验目的及对比方法

实验目的是考查本文提出的二类代价敏感纠错输出编码

框架BCS\_ECOC是否具有与现有方法相似或者更好的分类性能。实验将BCS\_ECOC在“一对一”编码矩阵和“一对多”编码矩阵下的代价敏感分类方法(分别简称为BCS-OVO, BCS-OVA)与已有同类方法进行比较。实验比较的方法包括CSOVO<sup>[17]</sup>, SECOC-OVO<sup>[16]</sup>, CSOVA<sup>[15]</sup>, SECOC-OVA<sup>[16]</sup>。文献[17]已经将CSOVO与其他方法进行了比较,这些方法包括Weighted All-pairs (WAP), Tree, Filter Tree (FT)<sup>[19]</sup>以及All-pair Filter Tree (APFT)<sup>[19]</sup>。文献[17]中,通过在大量代价敏感分类数据集上的实验,作者指出CSOVO的性能与WAP相当,但明显好于其他的几种方法。所以,在实验中,我们关注于BCS\_ECOC框架下的方法(BCS-OVO, BCS-OVA)是否拥有与CSOVO相似甚至更好的性能。此外,对于CSOVO, CSOVA, SECOC-OVO和SECOC-OVA<sup>1)</sup>,本文方法的源代码在此基础上进行了扩展。

### (2)实验分组设计

为了消除不同编码矩阵的影响,在每个数据集上将这些方法被划分为两组:一组是BC-CSOVO, CSOVO<sup>[17]</sup>和SECOC-OVO<sup>[16]</sup>;另一组是BC-CSOVA, CSOVA<sup>[15]</sup>和SECOC-OVA<sup>[16]</sup>。对每一组中的方法的性能进行相互比较。同时,引入了3种不同类型的多类代价矩阵<sup>[20]</sup>。对于每个数据集,只随机生成一个多类代价矩阵,以保证所有方法都在同一个多类代价矩阵下进行比较。

### (3)实验及参数

对于每一个数据集,利用10重交叉验证的结果来衡量每个方法的性能。采用10重交叉验证结果的误分类代价的平均值作为每个方法最终的误分类总代价。对于所有方法,二类代价敏感基分类器为代价敏感支持向量机,核函数为高斯径向基核函数,参数全部采用默认设置。需要指出的是,对于本文提出的方法,在人工数据集上以Bhattacharyya距离作为类间的可分性度量,在UCI数据集上以类内类间距离作为类间的可分性度量。同时,解码方法为最小汉明距离解码。

## 4.3 实验结果及分析

### (1)人工数据集

表3—表5分别列出了3种不同类型的多类代价矩阵<sup>[20]</sup>下,不同方法在人工数据集上的误分类总代价。

表3 不同方法在人工数据集上的误分类总代价(第一种类型多类代价矩阵)

Table 3 Test cost of different methods on artificial datasets using cost matrix in type a

Datasets	BCS-OVO	CSOVO	SECOC-OVO	BCS-OVA	CSOVA	SECOC-OVA
1#	3.00	2.70	5	4.50	3.60	5
2#	13.10	13.10	25	15.30	18.40	25
3#	25.00	25.00	50	26.80	31.60	50
4#	54.50	54.50	81.70	58.10	68.20	100
5#	85.60	87.60	174	99.30	122.30	150
6#	109.20	109.20	194	130.20	138.7	200
Average cost	48.40	48.68	88.28	55.70	63.80	88.33

表4 不同方法在人工数据集上的误分类总代价(第二种类型多类代价矩阵)

Table 4 Test cost of different methods on artificial datasets using cost matrix in type b

Datasets	BCS-OVO	CSOVO	SECOC-OVO	BCS-OVA	CSOVA	SECOC-OVA
1#	24.50	24.50	40.70	19.70	15.00	47
2#	100.50	100.50	301.30	106.70	124.30	265
3#	180.90	180.90	616	228.00	234.40	590
4#	216.10	215.70	807.80	264.60	248.40	740
5#	446.30	446.30	785.20	782.30	671.50	1200
6#	620.20	621.70	1173.10	780.60	758.70	1720
Average cost	264.75	264.93	620.68	363.65	342.05	760.33

<sup>1)</sup> <http://www.csie.ntu.edu.tw/~htlin/program/CSSVM>

表 5 不同方法在人工数据集上的误分类总代价(第三种类型多类代价矩阵)

Datasets	BCS-OVO	CSOVO	SECOC-OVO	BCS-OVA	CSOVA	SECOC-OVA
1#	31.80	<b>27.00</b>	67.40	38.60	<b>33.80</b>	83
2#	<b>102.60</b>	107.20	267.90	150.30	<b>104.50</b>	305
3#	286.20	<b>255.10</b>	533.10	453.20	<b>341.40</b>	600
4#	560	<b>447.90</b>	855	721.80	<b>563.20</b>	1100
5#	773.90	<b>543</b>	1541.50	939.70	<b>721.50</b>	1560
6#	863.60	<b>768.10</b>	2379.90	1327.20	<b>1115.40</b>	2720
Average cost	436.35	<b>358.05</b>	940.80	605.13	<b>479.96</b>	1061.30

表 3—表 5 中每个数据集上 3 种方法中的最小误分类总代价被加粗显示。每张表的最后一行给出了每个方法在所有数据集上的误分类总代价的平均值。

1)由表 3 和表 4 可以看出,BCS-OVO 和 BCS-OVA 拥有与 CSOVO 和 CSOVA 相似的误分类总代价的平均值。但在表 6 中,BCS-OVO 和 BCS-OVA 误分类总代价的平均值明显高于 CSOVO 和 CSOVA。即对于前两种类型的多类代价矩阵,本文方法的性能与 CSOVO 和 CSOVA 的性能相当甚至更好,但是在第三种类型的多类代价矩阵上性能较差。

2)由表 3—表 5 可以看出,本文提出方法的误分类总代价的平均值都显著低于 SECOC-OVO 和 SECOC-OVA 的平均值,性能明显更优。

#### (2)UCI 数据集

表 6—表 8 分别给出了 3 种不同类型的多类代价矩阵下,不同方法在 UCI 数据集上的误分类总代价。表中每个数据集上 3 种方法中的最小误分类总代价被加粗显示。每张表的最后一行给出了每个方法在所有数据集上的误分类总代价的平均值。

表 6 不同方法在 UCI 数据集上的误分类总代价(第一种类型多类代价矩阵)

Datasets	BCS-OVO	CSOVO	SECOC-OVO	BCS-OVA	CSOVA	SECOC-OVA
balance	<b>8.70</b>	11.90	57.70	<b>11.20</b>	11.30	57.60
glass	<b>3.40</b>	<b>3.40</b>	12.60	<b>4.20</b>	<b>4.20</b>	14.40
iris	4	10.50	6.10	<b>1.10</b>	<b>1.10</b>	10
letter	39	<b>33</b>	117.10	64.50	<b>59.40</b>	116.60
page-blocks	<b>27.80</b>	79.30	57.40	<b>32.80</b>	37.60	56
sat	118.60	<b>116.40</b>	489.80	123.40	<b>123.20</b>	490.20
segmentation	14.70	<b>13.20</b>	254.20	<b>175</b>	283.20	2640
soybean	17.50	<b>17.30</b>	30.20	<b>19</b>	<b>19</b>	29.60
thyroid	2.10	1	4.60	1.40	1.40	6.50
vehicle	32.90	<b>30.90</b>	79.50	30.40	<b>29.50</b>	64.70
vowel	<b>2.40</b>	2.50	98.90	4.70	<b>4.30</b>	90
wine	2.10	2.60	11.20	<b>3.30</b>	<b>3.30</b>	11.90
yeast	<b>62.50</b>	64.10	148.10	<b>66</b>	75.90	124
zoo	<b>2.70</b>	2.70	6	<b>2.90</b>	<b>2.90</b>	6
Average cost	<b>24.17</b>	27.77	98.10	<b>38.56</b>	46.87	265.53

表 7 不同方法在 UCI 数据集上的误分类总代价(第二种类型多类代价矩阵)

Datasets	BCS-OVO	CSOVO	SECOC-OVO	BCS-OVA	CSOVA	SECOC-OVA
balance	<b>31.90</b>	<b>31.90</b>	292.90	<b>30.70</b>	45.70	720
glass	<b>32.10</b>	<b>32.10</b>	197.50	<b>51.50</b>	51.60	205.30
iris	2.50	2.50	6.90	4.70	2.70	55
letter	350.40	<b>348.10</b>	1342.90	<b>705.50</b>	<b>705.50</b>	1342.90
page-blocks	<b>175.50</b>	301.60	3451.50	<b>152.70</b>	170.80	447
sat	<b>605.50</b>	<b>605.50</b>	5707.90	863.30	<b>729</b>	5825.10
segmentation	119.40	<b>108.10</b>	2541	192.40	<b>159.40</b>	2277
soybean	196.30	<b>190.50</b>	310.20	<b>205.30</b>	<b>205.30</b>	310.20
thyroid	2.50	2.50	22.10	2.80	<b>2.40</b>	30
vehicle	<b>165.80</b>	<b>165.80</b>	197.90	247.80	<b>187.40</b>	534.10
vowel	9	9	819	17.50	<b>17.50</b>	819
wine	<b>16.70</b>	<b>16.70</b>	148.20	<b>33.90</b>	<b>33.90</b>	132.60
yeast	569.20	<b>557.10</b>	940.50	653.90	<b>533.90</b>	940.50
zoo	21.50	<b>18.50</b>	83	<b>22.30</b>	<b>22.30</b>	39.80
Average cost	<b>164.16</b>	170.70	1147.30	227.45	<b>204.81</b>	977.03

表 8 不同方法在 UCI 数据集上的误分类总代价(第三种类型多类代价矩阵)

Datasets	BCS-OVO	CSOVO	SECOC-OVO	BCS-OVA	CSOVA	SECOC-OVA
balance	29.50	<b>20.40</b>	99.90	71.90	<b>40.60</b>	259.20
glass	29.40	<b>22.30</b>	105	<b>33.70</b>	34	105
iris	5.60	<b>4.80</b>	21.90	<b>6.40</b>	6.90	25
letter	309.10	<b>241.30</b>	1026.20	<b>541.50</b>	<b>541.50</b>	1026.20
page-blocks	498.30	<b>257.40</b>	746.50	<b>239.60</b>	299.80	735.70
sat	797.30	<b>366.40</b>	4003.10	763.10	<b>746.60</b>	4655.70
segmentation	<b>122.10</b>	150.10	1222.70	<b>254.30</b>	272.90	1914
soybean	196.40	<b>132.70</b>	321.80	<b>209.40</b>	<b>209.40</b>	321.80
thyroid	7	<b>6.50</b>	32.50	7.40	7.60	42.50
vehicle	277.10	<b>189.60</b>	767	333.70	<b>319.30</b>	991.40
vowel	6.30	<b>5.60</b>	864	<b>26.50</b>	<b>26.50</b>	864
wine	<b>17.30</b>	24.10	77.60	<b>34.40</b>	<b>34.40</b>	133.60
yeast	753.70	<b>623</b>	768.90	<b>646.20</b>	672.40	770.90
zoo	25.60	<b>16.30</b>	50.10	<b>25.10</b>	<b>25.10</b>	59.10
Average cost	219.62	<b>147.17</b>	721.94	<b>228.08</b>	231.21	850.29

1) 本文方法在 3 种类型的多类代价矩阵下取得了与 CSOVO 和 CSOVA 相似的结果。除了第三种类型的多类代价矩阵下, BCS-OVO 的误分类总代价的平均值要高于 CSOVO 的平均值, 这一结论与在人工数据集上的结论是一致的。

2) 本文方法与 SECOC-OVO 和 SECOC-OVA 相比性能明显更好。同时, 实验也验证了在实际问题中利用类内类间距离计算类间的可分性是可行的, 因为对于实际应用而言计算某一类的概率密度函数还是相当困难的。

通过以上实验可知: 本文提出的基于多类代价矩阵计算基分类器的二类代价矩阵的方法是有效的、可行的, 且基于此提出的新的二类代价敏感纠错输出编码 BCS\_ECOC 与现有方法相比具有相似甚至更好的分类性能。

**结束语** 本文研究了基于纠错输出编码实现多类代价敏感分类的方法, 提出一种新的代价敏感纠错输出编码框架; 将多类代价敏感分类问题分解为多个二类代价敏感分类问题, 为获得其中每个二类代价敏感基分类器的二类代价矩阵, 提出了利用已知多类代价矩阵计算每一类的误分类代价的期望值的方法, 并给出了为每个基分类器计算二类代价矩阵的通用的计算公式; 为验证提出方法的性能, 将其与现有方法在不同数据集上进行了比较, 实验结果表明, 提出方法具有相似甚至更好的性能。相比于 CSOVO 和 CSOVA 等特定的方法, BCS\_ECOC 作为一种通用框架, 可以应用现有的多种编码方法和解码方法, 从而为提高多类代价敏感分类的性能提供了可能。

## 参 考 文 献

- [1] ALI S, MAJID A, JAVED S G, et al. Can-CSC-GBE: Developing Cost-sensitive Classifier with Gentleboost Ensemble for breast cancer classification using protein amino acids and imbalanced data[J]. *Computers In Biology and Medicine*, 2016, 73: 38-46.
- [2] KIM Y J, BAIK B, CHO S. Detecting financial misstatements with fraud intention using multi-class cost-sensitive learning[J]. *Expert Systems With Applications*, 2016, 62: 32-43.
- [3] ZHOU Z H, LIU X Y. On multi-class cost-sensitive learning [J]. *Computational Intelligence*, 2010, 26(3): 232-257.
- [4] KULLUK S, ÖZBAKIR L, TAPKAN P Z, et al. Cost-sensitive meta-learning classifiers: MEPAR-miner and DIFACONN-miner [J]. *Knowledge-Based Systems*, 2016, 98: 148-161.
- [5] ZHANG G Q, SUN H J, et al. Cost-sensitive dictionary learning for face recognition[J]. *Pattern Recognition*, 2016, 60: 613-629.
- [6] JIANG L X. Cost-sensitive Bayesian network classifiers[J]. *Pattern Recognition Letters*, 2014, 45: 211-216.
- [7] BÁÑEZ A, BIELZA C, LARRANAGA P. Cost-sensitive selective naïve Bayes classifiers for predicting the increase of the h-index for scientific journals[J]. *Neurocomputing*, 2014, 135: 42-52.
- [8] CHEN Z, XIAO X Y, LI C S, et al. Real-time transient stability status prediction using cost-sensitive extreme learning machine [J]. *Neural Computing and Application*, 2016, 27: 321-331.
- [9] ZHANG L, ZHANG D. Evolutionary Cost-Sensitive Extreme Learning Machine[J]. *IEEE Transactions on Neural Networks and Learning Systems*, 2017, 28(12): 3045-3060.
- [10] DIETTERICH T, BAKIRI G. Solving multiclass learning problems via error-correcting output codes[J]. *Journal of Artificial Intelligence Research*, 1995, 34(2): 263-268.
- [11] LEI L, WANG X D. An Overview of Multi-Classification Based on Error-Correcting Output Codes[J]. *Acta Electronica Sinica*, 2014, 42(9): 1794-1800.
- [12] BAI X L, NIWAS S I, LIN W S, et al. Learning ECOC code matrix for multiclass classification with application to Glaucoma diagnosis[J]. *Journal of Medical Systems*, 2016, 40: 78.
- [13] LIU K H, ZENG Z H, NG V T Y. A hierarchical ensemble of ECOC for cancer classification based on multi-class microarray data[J]. *Information Sciences*, 2016, 349-350: 102-118.
- [14] SANTHANAM V, MORARIU V I, HARWOOD D, et al. A nonparametric approach to extending generic binary classifiers for multi-classification[J]. *Pattern Recognition*, 2016, 58(1): 149-158.
- [15] LIN H T. From ordinal ranking to binary classification[D]. California Institute of Technology, 2008.
- [16] LANGFORD J, BEYGEZIMER A. Sensitive error correcting output codes[J]. *Lecture Note in Artificial Intelligence*, 2005, 3559(1): 158-172.
- [17] LIN H T. Reduction from cost-sensitive multiclass classification to one-versus-one binary classification[C]// *JMLR: Workshop and Conference Proceedings*, 2014, 39: 371-386.
- [18] ASUNCION A, NEWMAN D. School of Information and Computer Sciences[M]// *UCI machine learning repository*. Irvine, CA, USA: University of California, 2007.
- [19] BEYGEZIMER A, LANGFORD J, RAVIKUMAR P. Error correcting tournaments[EB/OL]. <http://arxiv.org/abs/0902.3176>.
- [20] TING K M. An instance-weighting method to induce cost-sensitive trees[J]. *IEEE Transactions on Knowledge and Data Engineering*, 2002, 14(3): 659-665.
- [21] ZHANG L, ZHANG D. Evolutionary Cost-Sensitive Extreme Learning Machine[J]. *IEEE Transactions on Neural Networks and Learning Systems*, 2017, 28(12): 3045-3060.



**WU Chong-ming**, born in 1966, Ph. D., associate professor. His main research interests include machine learning and intelligent information processing.



**WANG Xiao-dan**, born in 1966, Ph. D., professor. Her research interests include machine learning, pattern recognition.